

PATENT

Attorney Docket No. 3359

PATENT APPLICATION

**METHODS AND COMPUTER SOFTWARE PRODUCTS FOR SELECTING
NUCLEIC ACID PROBES**

Inventors:

Rui Mei
A citizen of the United States of America
Residing at 234 Rodonovan Dr., Santa Clara, CA 95051

Teresa Webster
A citizen of the United States of America
Residing at 10002 Pescadero Creek Rd., Loma Mar, CA
94021

Assignee:

Affymetrix, Incorporated
A corporation organized under the laws of Delaware

Entity:

Large

Affymetrix, Inc.
Attn: Wei Zhou
3380 Central Expressway
Santa Clara, CA 95051
(408) 731-5000

**METHODS AND COMPUTER SOFTWARE
PRODUCTS FOR SELECTING NUCLEIC ACID PROBES**

5

RELATED APPLICATIONS

*Sub
P1* This application is related to U.S. Patent Application Number _____,
Attorney docket number 3367, entitled "Methods and Computer Software Products for
10 Predicting Nucleic Acid Hybridization Affinity", and U.S. Patent Application Number
_____, Attorney Docket Number 3373, entitled "Methods and Computer Software
Products for Selection Nucleic Acid Probes Using Dynamic Programming", filed
concurrently herewith. Both applications are incorporated herein by reference for all
purposes.

15

BACKGROUND OF THE INVENTION

The present invention relates to methods for designing nucleic acid probe arrays.
U.S. Patent No. 5,424,186 describes a pioneering technique for, among other things,
forming and using high density arrays of molecules such as oligonucleotides, RNA or
20 DNA), peptides, polysaccharides, and other materials. This patent is hereby incorporated
by reference for all purposes. There is still great need for methods, systems and software
for designing high density nucleic acid probe arrays.

SUMMARY OF THE INVENTION

In one aspect of the invention, computer implemented methods are provided for
25 selecting oligonucleotide probes. The methods including steps of a) predicting

hybridization intensities of a plurality of candidate probes, b) predicting quantitative responses of the candidate probes to the amount of their targets, c) selecting the probes from the candidate probes according to their hybridization intensities and quantitative response and d) spacing the probes along the sequence to avoid overlapping probes.

- 5 In some embodiments, the quantitative response is the slope of the response curve of a probe and hybridization intensity (I) is determined using the equation:

$$Ln(I) = \sum_{i=1}^{3N} W_i S_i + C_2$$

Or

$$Ln(I) = \sum_{i=1}^{3N} W_i S_i$$

wherein W_i is a weight coefficient; S_i is a functional of the sequence of a probe; N is the number of bases of a probe; and C_2 is a constant. In some embodiments, the weight coefficient is determined using multiple linear regression analysis.

In some preferred embodiments, the methods for selecting probes further include a step of predicting mismatch hybridization intensities of corresponding mismatch probes of the candidate probes and the selecting step is also based upon the mismatch

15 hybridization intensities. In some cases, the mismatch probes are different from their corresponding candidate probes in one base pair in the middle of their sequences. In preferred embodiments, the match hybridization intensities are predicted according to the sequences of the candidate genes. In some embodiments, mismatch hybridization intensities are determined according to the following equation:

$$Ln(I) = \sum_{i=1}^{3N} W_i' S_i + C_2' \quad \text{or}$$

$$Ln(I) = \sum_{i=1}^{3N} W_i' S_i$$

wherein said W_i' is a weight coefficient; S_i is a functional of the sequence of the perfect match probe; N is the number of bases of the probe; and C_2' is a constant, and I is the intensity of the mismatch probe.

The method of selecting probes may further include a step of calculating a unified quality score based upon predicted hybridization intensities.

In another aspect of the invention, computer software products are provided for selecting oligonucleotide probes. The software product includes computer program code for predicting hybridization intensities of a plurality of candidate probes; computer program code for predicting quantitative responses of the candidate probes to the amount of their targets; and computer program code for selecting said probes from said candidate probes according to said hybridization intensities and said quantitative response; and a computer readable media for storing said computer program codes.

In some embodiments, the quantitative response is the slope of the response curve of a probe. The hybridization intensity (I) may be determined using the equation:

$$Ln(I) = \sum_{i=1}^{3N} W_i S_i + C_2$$

Or

$$Ln(I) = \sum_{i=1}^{3N} W_i S_i$$

wherein said W_i is a weight coefficient; S_i is a functional of the sequence of a probe; N is the number of bases of a probe; and C_2 is a constant.

The weight coefficient is determined using multiple linear regression analysis.

The computer software product comprising computer program code for predicting

mismatch hybridization intensities of corresponding mismatch probes of said candidate probes and wherein said selecting step is also based upon said mismatch hybridization intensities. The method of Claim 13 wherein said mismatch probes are different from their corresponding candidate probes in one base pair in the middle of their sequence.

The mismatch hybridization intensities may be predicted according to the sequences of said candidate genes. In some embodiment, the mismatch hybridization intensities are determined according to the following equation:

$$Ln(I) = \sum_{i=1}^{3N} W'_i \cdot S_i + C'_2 \text{ or}$$

$$Ln(I) = \sum_{i=1}^{3N} W'_i \cdot S_i$$

wherein said W'_i is a weight coefficient; S_i is a functional of said sequence of said probe; N is the number of bases of said probe; and C'_2 is a constant. In one additional embodiment, the computer program code for selecting probes include computer program code for calculating a unified score for each probe.

In yet another aspect of the invention, a system for selecting nucleic acid probes is provided. The system include a processor; and a memory being coupled to the processor, the memory storing a plurality machine instructions that cause the processor to perform a plurality of logical steps when implemented by the processor, the logical steps including:

- a) predicting hybridization intensities of a plurality of candidate probes;
- b) predicting quantitative responses of the candidate probes to the amount of their targets;
- c) selecting the probes from the candidate probes according to the hybridization intensities and the quantitative response;
- d) spacing the probes along the sequence to avoid overlapping probes.

In some embodiments, the quantitative response is the slope of the response curve of the probe. The hybridization intensity (I) is determined using the equation:

$$\ln(I) = \sum_{i=1}^{3N} W_i S_i + C_2$$

Or

$$\ln(I) = \sum_{i=1}^{3N} W_i S_i$$

wherein said W_i is a weight coefficient; S_i is a functional of the sequence of a probe; N is the number of bases of a probe; and C_2 is a constant. The weight coefficient may be determined using multiple linear regression analysis.

In some preferred embodiments, the logic steps may further include predicting mismatch hybridization intensities of corresponding mismatch probes of the candidate probes and the selecting step is also based upon mismatch hybridization intensities. The mismatch probes are different from their corresponding candidate probes in one base pair in the middle of their sequences. The mismatch hybridization intensities may be predicted according to the sequences of said candidate genes. In some embodiments, the mismatch hybridization intensities are determined according to the following equation:

$$Ln(I) = \sum_{i=1}^{3N} W_i' S_i + C_2' \quad \text{or}$$

$$Ln(I) = \sum_{i=1}^{3N} W_i' S_i$$

wherein said W_i' is a weight coefficient; S_i is a functional of the sequence of the probe; N is the number of bases of the probe; and C_2' is a constant. The selecting step
 5 may also include a step of calculating a unified quality score based upon predicted hybridization intensities.

The present predictive methods are preferably used to select a collection of probes and an array upon which they are used.

BRIEF DESCRIPTION OF THE DRAWINGS

The accompanying drawings, which are incorporated in and form a part of this specification, illustrate embodiments of the invention and, together with the description, serve to explain the principles of the invention:

Figure 1 illustrates an example of a computer system that may be utilized to
 15 execute the software of an embodiment of the invention.

Figure 2 illustrates a system block diagram of the computer system of Figure 1.

Figure 3 illustrates a model system for probe sequence-based prediction of probe
 hybridization behavior (probe quality).

Figure 4 illustrates a basic physical model for probe target interaction.

Figure 5 shows an example of S_i for an exemplary probe.

Figures 6A and 6B show predicted relative ΔG for perfect match and mismatch probes.

Figure 7 shows an overall reaction of probe target formation.

Figure 8 shows concentration dependency of hybridization intensity.

5 Figures 9A and 9B show the relationship between probe-target binding affinity (K_{app}) and the slope (S).

Figure 10 shows an embodiment of a process for selecting probes.

Figure 11 shows a pool of candidate probes.

Figure 12 shows another embodiment of a process for selecting probes.

10 Figure 13 shows yet another embodiment of a process for selecting probes.

Figure 14 shows a process for obtaining weight coefficients.

Figure 15 shows yeast clones used to produce targets.

Figure 16 shows a Latin Square design.

Figure 17 shows Latin Square data sets from yeast_test_hyb chips.

15 Figure 18 shows a crossvalidation bootstrapping process.

Figures 19A, 19B, 20A and 20B show correlation between predicted and observed hybridization intensities for perfect match probes and mismatch probes.

Figure 21 shows hybridization intensity at different spike concentrations.

Figure 22 shows correlation between predicted and observed intensities over the entire concentration range.

Figure 23 shows predicted versus observed intensities for negative control target.

Figure 24 shows predicted versus observed slopes and improvement of correlation between the two slopes after filtering saturated probes.

Figure 25 shows prediction of hybridization of a human expression chip to human target sequences using weight coefficients generated from the yeast model system.

Figure 26 shows distribution of correlation coefficients.

Figure 27 shows selection of probes using dynamic programming.

Figure 28 compares different methods for selecting probes.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

Reference will now be made in detail to the preferred embodiments of the invention. While the invention will be described in conjunction with the preferred embodiments, it will be understood that they are not intended to limit the invention to these embodiments. On the contrary, the invention is intended to cover alternatives, modifications and equivalents, which may be included within the spirit and scope of the invention.

I. Glossary

“Nucleic acids,” according to the present invention, may include any polymer or oligomer of nucleosides or nucleotides (polynucleotides or oligonucleotides), which include pyrimidine and purine bases, preferably cytosine, thymine, and uracil, and adenine and guanine, respectively. See Albert L. Lehninger, *PRINCIPLES OF BIOCHEMISTRY*, at 793-800 (Worth Pub. 1982) and L. Stryer *BIOCHEMISTRY*, 4th

Ed., (March 1995), both incorporated by reference. Indeed, the present invention contemplates any deoxyribonucleotide, ribonucleotide or peptide nucleic acid component, and any chemical variants thereof, such as methylated, hydroxymethylated or glucosylated forms of these bases, and the like. The polymers or oligomers may be heterogeneous or homogeneous in composition, and may be isolated from naturally-occurring sources or may be artificially or synthetically produced. See U.S. patent application Serial No. 08/630,427 which is incorporated herein by reference in its entirety for all purposes. In addition, the nucleic acids may be DNA or RNA, or a mixture thereof, and may exist permanently or transitionally in single-stranded or double-stranded form, including homoduplex, heteroduplex, and hybrid states. Oligonucleotides and polynucleotides are included in this definition and relate to two or more nucleic acids in a polynucleotide.

“Probe,” as used herein, is defined as a nucleic acid, such as an oligonucleotide, capable of binding to a target nucleic acid of complementary sequence through one or more types of chemical bonds, usually through complementary base pairing, usually through hydrogen bond formation. As used herein, a probe may include natural (i.e. A, G, U, C, or T) or modified bases (7-deazaguanosine, inosine, etc.). In addition, the bases in probes may be joined by a linkage other than a phosphodiester bond, so long as the bond does not interfere with hybridization. Thus, probes may be peptide nucleic acids in which the constituent bases are joined by peptide bonds rather than phosphodiester linkages.

“Target nucleic acid” refers to a nucleic acid (often derived from a biological sample), to which the probe is designed to specifically hybridize. It is either the presence or absence of the target nucleic acid that is to be detected, or the amount of the target nucleic acid that is to be quantified. The target nucleic acid has a sequence that is complementary to the nucleic acid sequence of the corresponding probe directed to the target. The term target nucleic acid may refer to the specific subsequence of a larger nucleic acid to which the probe is directed or to the overall sequence (e.g., gene or mRNA) whose expression level it is desired to detect. The difference in usage will be apparent from context.

An “array” may comprise a solid support with peptide or nucleic acid probes attached to said support. Arrays typically comprise a plurality of different nucleic acids or peptide probes that are coupled to a surface of a substrate in different, known locations. These arrays, also described as “microarrays” or colloquially “chips” have been generally described in the art, for example, U.S. Pat. Nos. 5,143,854, 5,445,934, 5,744,305, 5,677,195, 6,040,193, 5,424,186 and Fodor et al., Science, 251:767-777 (1991). Each of which is incorporated by reference in its entirety for all purposes. These arrays may generally be produced using mechanical synthesis methods or light directed synthesis methods which incorporate a combination of photolithographic methods and solid phase synthesis methods. Techniques for the synthesis of these arrays using mechanical synthesis methods, such as ink jet, channel block, flow channel, and spotting methods which are described in, e.g., U.S. Pat. Nos. 5,384,261, and 6,040,193, which are incorporated herein by reference in their entirety for all purposes. Although a planar array

surface is preferred, the array may be fabricated on a surface of virtually any shape or even a multiplicity of surfaces. Arrays may be peptides or nucleic acids on beads, gels, polymeric surfaces, fibers such as fiber optics, glass or any other appropriate substrate, see U.S. Patent Nos. 5,744,305, 5,770,358, 5,789,162, 5,708,153, 6,040,193 and 5,800,992, which are hereby incorporated in their entirety for all purposes. Arrays may be packaged in such a manner as to allow for diagnostics or other manipulation of in an all inclusive device, see for example, US Patent Nos. 5,856,174 and 5,922,591, and 5,945,334, which are incorporated herein in their entirety by reference for all purposes. See also U.S. patent application Serial No. 09/545,207 which is incorporated herein in its entirety for all purposes for additional information concerning arrays, their manufacture, and their characteristics. It is hereby incorporated by reference in its entirety for all purposes.

II. Probe Selection Systems

As will be appreciated by one of skill in the art, the present invention may be embodied as a method, data processing system or program products. Accordingly, the present invention may take the form of data analysis systems, methods, analysis software, etc. Software written according to the present invention is to be stored in some form of computer readable medium, such as memory, or CD-ROM, or transmitted over a network, and executed by a processor. For a description of basic computer systems and computer networks, see, e.g., Introduction to Computing Systems: From Bits and Gates to C and Beyond by Yale N. Patt, Sanjay J. Patel, 1st edition (January 15, 2000) McGraw

Hill Text; ISBN: 0072376902; and Introduction to Client/Server Systems : A Practical Guide for Systems Professionals by Paul E. Renaud, 2nd edition (June 1996), John Wiley & Sons; ISBN: 0471133337.

Computer software products may be written in any of various suitable
5 programming languages, such as C, C++, Fortran and Java (Sun Microsystems). The computer software product may be an independent application with data input and data display modules. Alternatively, the computer software products may be classes that may be instantiated as distributed objects. The computer software products may also be component software such as Java Beans (Sun Microsystems), Enterprise Java Beans
10 (EJB), Microsoft® COM/DCOM, etc.

Figure 1 illustrates an example of a computer system that may be used to execute the software of an embodiment of the invention. Figure 1 shows a computer system 1 that includes a display 3, screen 5, cabinet 7, keyboard 9, and mouse 11. Mouse 11 may have one or more buttons for interacting with a graphic user interface. Cabinet 7 houses a
15 CD-ROM or DVD-ROM drive 13, system memory and a hard drive (*see* Figure 2) which may be utilized to store and retrieve software programs incorporating computer code that implements the invention, data for use with the invention and the like. Although a CD 15 is shown as an exemplary computer readable medium, other computer readable storage media including floppy disk, tape, flash memory, system memory, and hard drive may be
20 utilized. Additionally, a data signal embodied in a carrier wave (*e.g.*, in a network including the Internet) may be the computer readable storage medium.

Figure 2 shows a system block diagram of computer system 1 used to execute the software of an embodiment of the invention. As in Figure 1, computer system 1 includes monitor 3, keyboard 9, and mouse 11. Computer system 1 further includes subsystems such as a central processor 51, system memory 53, fixed storage 55 (*e.g.*, hard drive), removable storage 57 (*e.g.*, CD-ROM), display adapter 59, sound card 61, speakers 63, and network interface 65. Other computer systems suitable for use with the invention may include additional or fewer subsystems. For example, another computer system may include more than one processor 51 or a cache memory. Computer systems suitable for use with the invention may also be embedded in a measurement instrument.

III. Methods for Predicting Quality Scores of Probes

In a preferred embodiment, arrays of oligonucleotides or peptides, for example, are formed on the surface by sequentially removing a photoremovable group from a surface, coupling a monomer to the exposed region of the surface, and repeating the process. These techniques have been used to form extremely dense arrays of oligonucleotides, peptides, and other materials. The synthesis technology associated with this invention has come to be known as “VLSIPS™” or “Very Large Scale Immobilized Polymer Synthesis” technology and is further described below.

Additional techniques for forming and using such arrays are described in U.S. Patent Nos. 5,384,261, and 6,040,193 which are also incorporated by reference for all purposes. Such techniques include systems for mechanically protecting portions of a substrate (or chip), and selectively deprotecting/coupling materials to the substrate. Still

further techniques for array synthesis are provided in U.S. Application No. 08/327,512, also incorporated herein by reference for all purposes.

Nucleic acid probe arrays have found wide applications in gene expression monitoring, genotyping and mutation detection. For example, massive parallel gene expression monitoring methods using nucleic acid array technology have been developed to monitor the expression of a large number of genes (e.g., U.S. Patent Numbers 5,871,928, 5,800,992 and 6,040,138; de Saizieu *et al.*, 1998, Bacteria Transcript Imaging by Hybridization of total RNA to Oligonucleotide Arrays, NATURE BIOTECHNOLOGY, 16:45-48; Wodicka *et al.*, 1997, Genome-wide Expression Monitoring in *Saccharomyces cerevisiae*, NATURE BIOTECHNOLOGY 15:1359-1367; Lockhart *et al.*, 1996, Expression Monitoring by Hybridization to High Density Oligonucleotide Arrays. NATURE BIOTECHNOLOGY 14:1675-1680; Lander, 1999, Array of Hope, NATURE-GENETICS, 21(suppl.), at 3, all incorporated herein by reference for all purposes). Hybridization-based methodologies for high throughput mutational analysis using high-density oligonucleotide arrays (DNA chips) have been developed, see Hacia *et al.*, 1996, Detection of heterozygous mutations in BRCA1 using high density oligonucleotide arrays and two-color fluorescence analysis. Nat. Genet. 14:441-447, Hacia *et al.*, New approaches to BRCA1 mutation detection, Breast Disease 10:45-59 and Ramsey 1998, DNA chips: State-of-Art, Nat Biotechnol. 16:40-44, all incorporated herein by reference for all purposes). Oligonucleotide arrays have been used to screen for sequence variations in, for example, the *CFTR* gene (U.S. Patent Number 6,027,880, Cronin *et al.*, 1996, Cystic fibrosis mutation detection by hybridization to light-generated DNA probe

arrays. Hum. Mut. 7:244-255, both incorporated by reference in their entireties), the human immunodeficiency virus (HIV-1) reverse transcriptase and protease genes (U.S. Patent Number 5,862,242 and Kozal et al., 1996, Extensive polymorphisms observed in HIV-1 clade B protease gene using high density oligonucleotide arrays. Nature Med. 1:735-759, both incorporated herein by reference for all purposes), the mitochondrial genome (Chee et al., 1996, Accessing genetic information with high density DNA arrays. Science 274:610-614) and the BRCA1 gene (U.S. Patent Number 6,013,449, incorporated herein by reference for all purposes).

Sub 10
In one aspect of the invention, a physical model that is based on the thermodynamic properties of the sequence is used to predict the array-based hybridization intensities of the sequence. Hybridization propensities may be described by energetic parameters derived from the probe sequence, and variations in hybridization and chip manufacturing conditions will result in changes in these parameters that can be detected and corrected. U.S. Patent Application Number _____, docket Number 3367, filed concurrently herewith and incorporated herein by reference, discloses methods for predicting nucleic acid hybridization affinity.

20 The values of weight coefficients in the physical model may be determined by empirical data because these values are influenced by assay conditions, which include hybridization and target fragmentation, and probe synthesis conditions, which include choice of substrates, coupling efficiency, etc.

In one embodiment (Figure 3), a model experimental system is used to generate empirical data and a computational model is used to process these data to solve for the

weight coefficients of the physical model. These solved weight coefficients are in turn placed back into the physical model, enabling it to predict the hybridization behaviors of new sequences.

The interaction between a probe and its target is described in Figure 4.

- 5 Basically, a target (T) hybridizes to its complementary probe (P) to form a probe-target duplex (P•T) (Figure 4), and the reaction is accompanied with favorable free energy change (Figure 4). The amplitude of the free energy change (ΔG) determines the stability of probe-target duplex. The duplex stability can be described by equilibrium constant (K_s), which is sequence-dependent. The relationship between K_s and ΔG may be given
- 10 by Boltzmann's equation:

$$K_s = \frac{k_{on}}{k_{off}} = e^{-\Delta G / RT} \quad [\text{Equation 1}]$$

where k_{on} and k_{off} are the rate constants for association and dissociation, respectively, of the probe-target duplex, R is the gas constant and T is the absolute temperature.

According to Equation 1, ΔG is a function of the sequence. The dependence of ΔG on

- 15 probe sequence can be quite complicated, but relatively simple models for ΔG have yielded good results.

Sup
P3

There are a number of ways to establish the relationship between the sequence and ΔG . In preferred embodiments, one model (equation 2), shown in U.S. Application Serial Number _____, Attorney Docket Number 3367, previously incorporated by

20 reference is shown below:

$$\Delta G_{seq} = \sum_{i=1}^{3N} P_i S_i \quad [\text{Equation 2}]$$

or

$$\Delta G_{seq} = \sum_{i=1}^{3N} P_i S_i + C \quad [\text{Equation 3}]$$

where N is the length (number of bases) of a probe. P_i is the value of the i th parameter
5 which reflects the ΔG of a base in a given sequence position relative to a reference base in
the same position. In preferred embodiments, the reference base is A. In this case, the
 P_i 's will be the free energy of a base in a given position relative to base A in the same
position. Figure 5 shows an example of how the value of S_i is determined based upon the
sequence of a probe. In this example, a probe, GTCA has $N=4$ and thus, it has $3 \times 4 = 12$ S_i
10 values. Each probe base position has three S values, each for a different possible base. In
this example, possible bases are evaluated in the sequence of C, G, and T (A is the
reference base). However, one of skill in the art would appreciate that the assignment of
this particular of base sequence is arbitrary. Alternatively evaluation sequence, such as
G, C, and T may also be used as long as the same scheme is used for model building and
15 for hybridization affinity prediction.

Based on the simple hybridization scheme described in Figure 4, the hybridization
intensity is proportional to the concentration of probe-target duplex, where C_0 is constant
(Equation 4). Under equilibrium condition, the intensity is directly related to ΔG
(Equation 5). This relationship is also expressed in natural logarithm form, where C_1 and
20 C_2 are constants (Equation 7) and Equation 6 also holds for approaching equilibrium
cases. According to Equation 2, the relationship between intensity and probe sequence is
described in Equation 7 and 8:

$$I = C_0 [P \cdot T] \quad [\text{Equation 4}]$$

$$[P \cdot T] = K_s [P][T] = e^{-\Delta G/RT} [P][T] \quad [\text{Equation 5}]$$

$$\text{Ln} I = -\Delta G/RT + \text{Ln}\{C_0[P][T]\} \quad [\text{Equation 6}]$$

$$\text{Ln} I = C_1 \sum_{i=1}^{3N} P_i S_i + C_2, \text{ where } C_2 =$$

$$\text{Ln}\{C_0[P][T]\} \text{ and } C_1 = -1/RT \quad [\text{Equation 7}]$$

$$\text{Ln} I = \sum_{i=1}^{3N} C_1 P_i S_i + C_2 = \sum_{i=1}^{3N} W_i S_i + C_2 \quad [\text{Equation 8}]$$

where $W_i = C_1 P_i$. The following is a linear regression model for probes of N bases in length using a training data set that contains intensity values of M probes.

$$\text{Ln}(I_1) = W_1 S_{11} + W_2 S_{21} + \dots W_{3N} S_{3N1}$$

$$\text{Ln}(I_2) = W_1 S_{12} + W_2 S_{22} + \dots W_{3N} S_{3N2}$$

...

$$\text{Ln}(I_M) = W_1 S_{1M} + W_2 S_{2M} + \dots W_{3N} S_{3NM}$$

Hybridization intensities (relative to a reference base, such as an A) for each type of bases can be solved at each position in the probe sequence may be predicted. Multiple linear regression analysis is well known in the art, see, for example, the electronic statistic book (<http://www.statsoftinc.com/textbook/stathome.html>); Darlington, R. B. (1990).

15 *Regression and linear models*. New York: McGraw-Hill, both incorporated by reference for all purposes. Computer software packages, such as SAS, SPSS, and MatLib 5.3

provide multiple linear regression functions. In addition, computer software code examples suitable for performing multiple linear regression analysis are provided in, for example, the Numerical Recipes (NR) books developed by Numerical Recipes Software and published by Cambridge University Press (CUP, with U.K. and U.S. web sites).

5 In a preferred embodiment, a set of probes of different sequences (probes 1 to M) is used as probes in experiments(s). Hybridization affinities (relative ΔG or $\ln(I)$) of the probes with their target are experimentally measured to obtain a training data set (see, example section *infra*). Multiple linear regression may be performed using hybridization affinities as $I [I_1 \dots I_m]$ to obtain a set of weight coefficients: $[W_1 \dots W_N]$. The weight
10 coefficients are then used to predict the hybridization affinities using Equation 7. Figure 6A shows relative predicted ΔG at every base position in a probe of 25 bases in an exemplary experiments (see the example section *infra* for a detailed description of experimental conditions).

In addition, in some embodiments, by using intensities derived from mismatch
15 probes that are probes designed to contain one or more mismatch bases from a reference probe, a set of weight coefficients may be obtained to predict the mismatch intensity using perfect match probe sequence. Figure 6B shows an example for predicting mismatch hybridization affinity at center base position.

Since other interactions such as probe self-folding, probe-to-probe interaction,
20 target self-folding and target-to-target interaction also interfere with the probe-target duplex formation, their contributions to the values of the weight coefficients may also be considered. Figure 7 shows an overall equilibrium scheme including the formation of a

probe-target duplex (PT), probe self-folding (P_F) and probe dimerization (PP). Probe folding renders the probe unavailable for binding with the target. Probe dimerization renders two probes unavailable for binding with the target. In some embodiments, the hybridization affinity prediction model accounts for probe folding and probe

5 dimerization:

$$\Delta G_{overall}^0 = -W_d \Delta G_d^0 + W_{PF} \Delta G_{PF}^0 + W_{PP} \Delta G_{PP}^0 \quad [\text{Equation 9}]$$

$$\ln I = C_1 \Delta G_{overall}^0 + C_2 \quad [\text{Equation 10}]$$

where W_d is the weight for sequence based probe affinity; W_{PF} is the weight for probe formation and W_{PP} is the weight for probe dimerization. Any methods that are capable of predicting probe folding and/or probe dimerization are suitable for at least some embodiments of the invention for predicting the hybridization intensity in at least some embodiments of the invention. In a particularly preferred embodiment, Oligowalk (available at <http://rna.chem.rochester.edu/RNAstructure.html>, last visited Nov. 3, 2000) may be used to predict probe folding.

One important criterion of probe selection for a quantitative gene expression assay is that hybridization intensities of the selected probes must correspond to target concentration changes. In some embodiments, the relationship between concentrations and intensities of a probe is modeled as:

$$\ln(I) = S \ln C + \ln K_{app} \quad [\text{Equation 11}]$$

20 or

$$I = K_{app} C^S \quad [\text{Equation 12}]$$

where I is intensity; K_{app} is apparent affinity constant; C is concentration of the target; and S is an empirical value corresponding to the slope of the line relates $\ln I$ and $\ln C$ ($0 < S < 1$) (see Figure 8).

Equation 12 describes the relationship between hybridization intensities of probes

5 and target concentration. For example, when S is equal to 1, the intensities of a probe linearly correspond to its target concentration (Figure 8). Thus, based on the S values of the probe, one can select probes that have good concentration dependence. Figure 9A shows the polynomial relationship between S and $\ln K_{app}$, indicating that when the value of $\ln K_{app}$ increases to a certain level the value of S reaches a plateau before starting to
10 decrease. This relationship allows the identification of not only low hybridization affinity probes (Figure 9B, bottom lines) but also GC-rich probes that have high affinity but bind to both specific and non-specific targets (Figure 9B, top line). These GC-rich probes have high intensities, but the intensities maintain constant when target concentration changes (Figure 9B, top line). Therefore, these probes have small slopes. In some
15 embodiments, linear regression modeling alone will not identify probes with a high propensity to saturate. That is because the linear model for each target concentration will predict the intensity that a probe would have had if it could bind to unlimited amount of target. Therefore, the predicted slope can be quite high when the observed slope is low (Figure 24, top). The well-behaved relationship between predicted $\ln K_{app}$ and observed
20 slope allows filtering probes with a high propensity to saturate based on the predicted $\ln K_{app}$ for the given probe. If $\ln K_{app}$ is above a cutoff value (e.g., 5, 6, 7, 8, 9, or 10, Figure 9), then the probe is effectively filtered as a candidate for probe selection. Figure

24 (middle) shows the predicted slope profiles after filtering as well as the significant improvement in the overall correlation after these regions are removed.

IV. Methods and Software for Selecting Probes

5 Figure 10 shows a computer-implemented process for selecting probe sequences from a pool of candidate probes. In this particular embodiment, the sequences of a pool of candidate oligonucleotide probe are processed by a quality predictor (101).

Sub
Fig 6

10 Throughout this application, the term probe may refer to the sequence of a probe. The pool of candidate oligonucleotide probes may be all possible probes against a particular target or targets. Typically, oligonucleotide probes are at least 10, 15, 20, 25 and 30 bases in length. Polynucleotide probes can be more than 10, 20, 25, 30, 100, 200, 500, 1000, or 5000 bases in length. Figure 11 illustrates a complete pool of candidate oligonucleotide probes (unfiled rectangle boxes) against a target (black rectangle box). Each of the probes is designed to be complementary to the target sequence. In this particular

15 embodiment, the oligonucleotides are 25mers. The first probe is complementary to bases 1-25 (from the 5' end) of the target sequence. The second probe is complementary to bases 2-26 and so on. While a complete pool is often desirable, it is not necessary to have a complete pool for at least some embodiments of the invention. In some cases, filters may be used to remove some of the probes from the pool.

20 The input to the quality predictor (Figure 10, 101) is the sequence of a pool of candidate probes. One of skill in the art would appreciate that the format of input is not critical. In some embodiments, the probe sequences may be inputted from one or a

number of probe sequence files. The file(s) may be plain text file(s), in the FASTA format or other suitable file format. Alternatively, the input may be a stream from other sources such as a data pocket stream from a remote networked computer.

5 ~~Sup A7~~ The quality predictor is a software module that calculates quality scores (the term score refers to any qualitative and quantitative values with regard to desired properties of a probe) for probes based upon the sequences of probes. In some embodiments, the quality score may include predicted values such as perfect match intensity, mismatch intensity and/or slope.

10 Probe selection module (103) selects probes based upon their scores. In preferred embodiments, the quality scores are combined to obtain a unified score. In some cases, the unified quality score is the simple summation of quality scores (e.g., Unified Quality Score=Perfect Match Intensity+Mismatch Intensity+Slope). The selection of probes may be based upon the scores only. For example, if certain number of probes are desired, the probes with the highest scores are selected until enough number of probes are selected.

15 Alternatively, a threshold-unified score may be established. Probes that have scores higher than the threshold score are selected.

20 ~~Sup A10~~ In preferred embodiment, the goal of probe selection step is to find the best probes to represent a sequence. The probe selection software module takes a set of probes and a set of quality measures for each probe. It then implements an optimization algorithm to find the best n probes, spread out across the gene. Methods for probe selection using optimization algorithm is described in U.S. Application Number _____, Docket

Number 3373, filed concurrently herewith and incorporated herein by reference in its entirety for all purposes.

5 *Sub 101* Figure 12 shows another embodiment of the computer implemented probe selection process of the invention, target sequences are inputted to a candidate probe generator (121) which produce either all possible probes of certain length or a subset of the all possible probes. The candidate probe sequences are fed to the quality score predictor (122) for calculating quality measures (scores, e.g., perfect match intensity, mismatch intensity and/or slope). The candidate probe sequences are also fed to a 3' bias score predictor (123) to obtain 3' bias scores that indicates the distance of probe sequence from the 3' end of target sequence. Since the current target preparation method is 3' biased, it is important to select probes that fall into range where its target will be made. The probe sequences may optionally be inputted into a cross hybridization score predictor (124) to calculate cross hybridization scores. The quality scores, 3' bias scores and/or cross hybridization score are combined by a probe score calculator module (125) to produce a unified score. A probe selection module (126) picks the probes with the score which indicates that _____.

20 Figure 13 shows a complete computer implemented probe selection process. In this preferred embodiment, target sequences (131) are used to generated a pool of candidate probes. The probe sequences are stored in a FASTA sequence file. A sequence file splitter (132) divides probe sequences to .seq file which store one sequence per file. The .seq files are processed using a OligoWalker batch tool (133) to produce a .rep file, one for each probe sequence. The .rep files contain ΔG values for the probes. The rep

files are inputted into a quality predictor (134). The quality predictor is based upon a multiple linear regression models derived from experiment data using, for example, yeast test chips (see also, example section below) (1310). The quality predictor calculates quality scores (measures, perfect match intensity, mismatch intensity and slope) as

5 described above in section II. The rep file is also inputted into a 3' bias score predictor (135) to estimate 3' bias scores for the probes.

SB
A12)

The multiple probe FASTA sequence file is also inputted into a cross hybridization predictor (136) to predict a cross hybridization score. The cross hybridization score predictor is based upon models (such as multiple linear regression

10 models) derived from experiment data (1311). In some embodiments, cross hybridization may also be evaluated by pruning probe sequences against a human genome data base (1312) which may be residing locally, in a local area network or in a remote site such as the Genbank (<http://www.ncbi.nlm.nih.gov>).

The quality measures, 3' bias scores and cross hybridization scores are combined

15 by the probe score calculator (137) to produce a unified score for each probe. The combined score is then used for selecting probes (138). The probe selection module takes a set of probes and a set of quality measures for each probe. It then implements a dynamic programming algorithm to find the best n probes, spread out across the gene. The selected probe sequences are stored in .101 files (139).

20 The following tables describe the various software modules in the exemplary embodiments described in Figure 13.

1: Multiple linear regression modeling tool

Description	Calculates the weights for the regression model. Its is a one time calculation. The results of the calculations will be used every time a new chip is designed.
Input	Yeast Test Chip, available from Affymetrix, Santa Clara, CA
Output	Multiple linear regression models, a set of weights.
Part of chip design	In this embodiment, it is not part of the software package for chip design. It is used as one time external process. However, in other exemplary embodiments, it may also become part of the software.

2: Sequence file splitter

Description	Splits a FASTA file of sequences into several sequence files one for each sequence in the instruction file. If max files in folder is greater than 0, subfolders are created in the output path. Each subfolder get up to the maximum files specified.
Input	<ul style="list-style-type: none">• FASTA file• Instructions file• Output path• Max files in one folder
Language / Tool	Java

3: Oligo Walk batch tool

Description	Runs Oligo Walk in batch mode. Oligo Walk produces a .rep file for each sequence. The .rep file contains a delta G values for each probe
Input	Batch of .seq files
Output	.rep file. The .rep file identifies a probe by a number and a sequence. The sequence is a reverse complement of the 25-mer it represents on the input sequence. The number is the beginning of the probe.
Part of chip design	Yes
Language / Tool	Microsoft ® Visual Basic

4: Quality predictor

Description	Takes in the MLR model measures and delta G values from Oligo Walk and produces 3 quality measures, perfect match intensity, mismatch intensity and slope.
Input	.rep file produced from Oligo Walk
Output	3 Quality measures for each probe. The probe is described as in the input format.
Part of chip design	Yes
Language / Tool	C

5: Cross Hyb Modeling Tool

Description	Analyzes the results of the yeast cross hyb chip to create a model for predicting the cross hyb score for a probe, based on the number of mismatches and positions of mismatches with 1 or more matching sequences.
Input	Results from the cross hyb chip
Output	A model that relates number of mismatches and positions of mismatches to a cross hyb score.
Part of chip design	In some embodiments, it is not part of the chip design package. Alternatively, it can be part of the package.

6: Cross Hyb Score Predictor

Description	Predicts a cross hyb score for a given set of probes. Its does so by matching the given probes with a genome and assigns a numeric score using the cross hyb models.
Input	<ul style="list-style-type: none"> ▪ Cross hyb models ▪ A genome ▪ Set of probes
Output	List of probes and corresponding cross hyb scores
Part of chip	Yes

design	
---------------	--

7: 3' Bias Score Predictor

Description	<p>Predicts the 3' bias score for a given set of probes. Earlier it was believed that most sequences have a sigmoid graph for the 3' bias.</p> <p>But, recently used sequences do not always follow the pattern</p> <p>Therefore, it is important to first study the 3' bias effect and then design a measurement model.</p>
Input	<ul style="list-style-type: none"> ▪ Set of probes
Output	List of probes and corresponding 3' bias scores
Part of chip design	Yes

8: Probe Score Calculator

Description	<p>Given a set files with probe information and scores, this program matches each probe in each sequence and calculates 1 unified score for each probe.</p>
Input	<ul style="list-style-type: none"> ▪ Set of probes ▪ Set of measures for each probe, each in a different file(s) <ul style="list-style-type: none"> ▪ 3 quality scores (probes defined in OligoWalk format) ▪ cross hyb score (probes defined in chip design format) ▪ 3' bias score (probes defined in chip design format)

Output	List of probes with a corresponding score.
Part of chip design	Yes

9: Probe selection algorithm

Description	Finds the best probes to represent a sequence. It takes a set of probes and a set of quality measures for each probe. It then implements a dynamic programming algorithm to find the best n probes, spread out across the gene.
Input	<ul style="list-style-type: none"> ▪ Set of probes ▪ Set of measures for each probe <ul style="list-style-type: none"> ▪ 3 quality scores ▪ cross hyb score ▪ 3' bias score ▪ Number of probes to choose
Output	.llq file
Part of chip design	Yes
Language / Tool	C

10: Algorithm Test Tool

Description	Tests the new probe selection algorithm. The probe selection algorithm is used to select probes for the known, Yeast test chip. The selected probes are analyzed for their intensity, slope and discrimination values on the yeast test chip.
Input	<ul style="list-style-type: none">▪ Probes selected for the sequences on the yeast test chip▪ Results from the yeast test chip

V. Examples

The following examples demonstrate the effectiveness of the methods of the invention for predicting hybridization intensities and for selecting oligonucleotide probes for gene expression monitoring.

A. Example 1: Prediction of Hybridization Intensities of Probes Against Yeast Genes

Figure 14 shows the overall process of the experiments. Yeast was used as a model system for this experiment because the yeast genome had been sequenced. Arrays containing nucleic acid probes complementary to yeast genes are commercially available from Affymetrix (Santa Clara, California), and include _____.

Genes were selected to cover sequence complexity such as GC content, secondary structure, Motif and gene clusters. Twenty probe pairs (perfect match and mismatch probes) were selected to cover entire sequence of one of the 112 selected yeast genes.

The probes are synthesized in situ on glass substrate using photo-directed synthesis

method that was disclosed in, for example, U.S. Patent Nos. 5,384,261, and 5,744,305, 5,445,934 and 6,040,138.

One hundred and twelve yeast clones representing the 112 genes were randomly divided into 14 groups (Figure 15). Labeled targets prepared from these clones were used as spikes for 14 experiments at various concentration levels from 0pM to 1024pM. In some experiments, the spikes derived from yeast gene clones were combined with labeled nucleic acid representing human complex background. A 14 x 14 Latin square design (Figure 16) was employed. The numbers in the table indicates the concentration used (pM). For each experiment, 14 groups of genes at 14 different concentrations were pooled together and hybridized to an oligonucleotide probe array. For each Latin Square 14 oligonucleotide probe array hybridization experiments were performed. Figure 17 shows experiments conducted.

Cross-validation (Figure 18) was used to evaluate the prediction. The cross-validation process held one gene for test and used the other 111 genes to solve the weight coefficients that in turn were used to predict intensities for the test genes, as described in Figure 14. The correlation between the predicted and measured intensity for one test gene (YDR113C) is shown in Figure 19A and Figure 19B shows the correlation against target sequence, where lines represent the predicted values and dots represent the observed values. The correlation of the predicted and measured values for perfect match (PM) and mismatch (MM) probes is also demonstrated in Figure 20A and 20B respectively, where lines represent the predicted values and dots represent the observed values for gene YGR109C.

Figure 21 shows predicted intensity versus actual intensity at various target spike concentrations, where lines indicate the predicted values and dots represent observed values. Figure 22 shows correlation coefficients between predicted and observed intensity (LnI) as function of concentration, where top and bottom lines represent perfect matches and mismatches, respectively. The high correlation (0.85) holds for 4000-fold concentration range (Figure 22), and the results demonstrate that the methods of invention are able to predict probe behaviors through a wide dynamic range.

Figure 23 shows predicted versus observed intensities when the target transcripts were derived from genes in the wrong orientation, which results no complimentary target was generated for the probes. As shown in Figure 24, predicted intensities (lines) had no correlation with observed intensity (dots) because right target is absent. The result indicates the prediction method is accurate and specific.

Figure 24 shows predicted slope versus observed slope. In some regions in the Figure 24(top), the values of predicted slope (lines) can be quite high when the values of observed slope (dots) because of the saturated probes in those regions. According to Equation 12 and Figure 9, the saturated probes can be identified and removed. Figure 24 (middle) shows the predicted slope profiles after filtering the saturated probes and the significant improvement in the overall correlation after these regions are removed.

B. Example 2: Prediction of Hybridization Intensities of Probes from Human Genes

This example demonstrates that weight coefficients obtained from the model yeast experiment system is also able to predict the intensities on the human gene expression

chip and the predicted intensities (left bar) are highly correlated with observed intensities (right bar) at each probe position as indicated by x-axis. The correlation is shown in Figures 25 A-E. Typically, the correlation coefficients ranged from 0.45-0.83. The distribution of the correlation coefficients are shown in Figure 26. These results demonstrate that the probe selection model may be generalized to different organisms such as mammals, plants,

C. Example 3: Probe Selection

This example demonstrates that the model-based probe selection method and software may provide improvement over current probe selection methods. Figure 27 shows intensity values of sixteen probes (open squares) selected for the Yer161c gene based upon quality scores and using dynamic programming. Figure 27 also shows that the sixteen selected probes (open squares) are spaced along sequence. Figure 28 shows a comparison of average intensity difference (between perfect match and mismatch) values of probe selected by various methods for all yeast test genes. Probes selected randomly (diamonds) were similar to those selected according empirical rules (squares). The model based selection method (triangles) improved average intensity difference values. The result indicates the model-selected probes have high sensitivity and specificity.

Conclusion

The present invention provides methods and computer software products for predicting nucleic acid hybridization affinity, detecting mutation, selecting better-behaved

probes, and improving probe array manufacturing quality control. It is to be understood that the above description is intended to be illustrative and not restrictive. Many variations of the invention will be apparent to those of skill in the art upon reviewing the above description. By way of example, the invention has been described primarily with reference to the use of a high density oligonucleotide array, but it will be readily recognized by those of skill in the art that the methods may be used to predict the hybridization affinity of other immobilized probes, such as probes that are immobilized in or on optical fibers or other supports by any deposition methods. The basic methods and computer software of the invention may also be used to predict solution-based hybridization. The scope of the invention should, therefore, be determined not with reference to the above description, but should instead be determined with reference to the appended claims, along with the full scope of equivalents to which such claims are entitled.

All references cited herein are incorporated herewith by reference for all purposes.